# Prediction Serving

*what happens after learning?*

**Joseph E. Gonzalez**

Asst. Professor, UC Berkeley

jegonzal@cs.berkeley.edu

Co-founder, GraphLab (now Turi Inc.)

joseph@dato.com

Prediction Serving

Graph Systems

Graph Frames

Learning Systems

Time Series

Frequency Domain Analytics Systems

Cluster Management

Multi Task Learning for Job Scheduling

Cross-Cloud Perf. Estimation

# Outline



**Active Collaborators**

Daniel Crankshaw

Xin Wang

Michael Franklin

Ion Stoica

# Learning



**Timescale:** minutes to days
**Systems:** offline and batch optimized
*Heavily studied ... major focus of the **AMPLab***

# Learning

# Inference

Big Data

Training

Big Model

Query

Decision

Application

**Timescale:** ~10 milliseconds
**Systems:** *online* and *latency* optimized
*Less studied ...*

# Learning

# Inference

Big Data

Training

Big Model

Query

Decision

Application

Feedback

**Learning** | **Inference**

Training

Decision

Big Data

Application

**Timescale:** hours to weeks
**Systems:** combination of systems
Less studied …

Feedback

# VELOX Model Serving System [CIDR'15]

*Daniel Crankshaw, Peter Bailis, Haoyuan Li, Zhao Zhang, Joseph Gonzalez, Michael J. Franklin, Ali Ghodsi, and Michael I. Jordan*

**Adaptive
(~1 seconds)**

**Responsive
(~10ms)**

# Key Insight:

*Decompose models into fast and slow changing components*

# Hybrid Offline + Online Learning

Update feature functions offline using batch solvers
- Leverage high-throughput systems (Tensor Flow)
- Exploit slow change in population statistics

$$f(x; \theta)^T w_u$$

Update the user weights online:
- Simple to train + more robust model
- Address rapidly changing user statistics

# Common modeling structure

$$f(x;\theta)^T w_u$$

Matrix Factorization



Deep Learning



Input

Ensemble Methods

# Velox Online Learning for Recommendations
## (Simulated News Rec.)



**Partial Updates:** *0.4 ms*
**Retraining:** *7.1 seconds*

*>4 orders-of-magnitude*
***faster adaptation***

# ⚛ *VELOX*: the Missing Piece of BDAS



Learning

amplab

| Spark Streaming | BlinkDB | Graph Frames | Keystone ML |
| | Spark SQL | GraphX | MLLib |

Spark

Mesos

Tachyon

HDFS, S3, …

**B** erkeley
**D** ata
**A** nalytics
**S** tack

# VELOX: the Missing Piece of BDAS

# VELOX Architecture

Fraud Detection

Content Rec.





**Keystone ML**

**MLLib**

**Spark**

*Velox*

Single JVM Instance

# VELOX Architecture

| Fraud Detection | Content Rec. | Personal Asst. | Robotic Control | Machine Translation |
|---|---|---|---|---|



**Velox**

- Keystone ML
- MLLib
- Spark

Single JVM Instance

scikit learn

Dato Create

Caffe

dmlc mxnet

TensorFlow

VW

KALDI

theano

# Clipper **Generalizes** Velox Across ML Frameworks

Fraud
Detection

Content
Rec.

Personal
Asst.

Robotic
Control

Machine
Translation



## Clipper

# Clipper

**Key Insight:**

*The challenges of prediction serving can be addressed between end-user applications and machine learning frameworks*

As a result, Clipper is able to:

➢ **hide complexity**

  ➢ by providing a *common prediction interface*

➢ **bound latency** and **maximize throughput**

  ➢ through *approximate caching* and *adaptive batching*

➢ enable *robust* **online learning** and **personalization**

  ➢ through generalized *split-model correction policies*

**without modifying** *machine learning frameworks or end-user applications*

# Clipper Design Goals

Low and **bounded** latency predictions

> ➤ interactive applications need reliable latency objectives

Up-to-date and personalized predictions **across models** and **frameworks**

> ➤ generalize the split model decomposition

Optimize **throughput** for performance under heavy load

> ➤ single query can trigger many predictions

**Simplify** deployment

> ➤ serve models using the original code and systems

# Clipper Architecture

Fraud
Detection

Content
Rec.

Personal
Asst.

Robotic
Control

Machine
Translation



## Clipper

# Clipper Architecture



Applications

**Predict** ↕ **RPC/REST Interface** ↕ Observe

## Clipper

# Clipper Architecture

# Clipper Architecture

Applications

Predict ↕  **RPC/REST Interface**  ↕ Observe

## Clipper

*Improve accuracy through **ensembles**, **online learning** and **personalization***

### Correction Layer

*Provide a **common interface** to models while **bounding latency** and **maximizing throughput**.*

### Model Abstraction Layer

RPC ↕    RPC ↕    RPC ↕    RPC ↕

Model Wrapper (MW)    MW    MW    MW

● ● ●

# Clipper Architecture

Applications

Predict ↕ **RPC/REST Interface** ↕ Observe

## Clipper

**Correction Layer**

Correction Policy

**Model Abstraction Layer**

Approximate Caching

Adaptive Batching

RPC ↕     RPC ↕     RPC ↕     RPC ↕

Model Wrapper (MW)    MW    MW    MW    ●●●

Model Abstraction Layer

Approximate Caching

Adaptive Batching

RPC — Model Wrapper (MW) — KeystoneML

RPC — MW — Caffe

RPC — MW — TF

RPC — MW — scikit learn

• • •

Provides a unified generic prediction API across **frameworks**

➤ **Reduce Latency** → Approximate Caching

➤ **Increase Throughput** → Adaptive Batching

➤ **Simplify Deployment** → RPC + Model Wrapper

**Model Abstraction Layer**

Approximate Caching

Adaptive Batching

RPC

RPC

RPC

RPC

Model Wrapper (MW)

MW

MW

MW

KeystoneML

Caffe

●●●

**Common Interface → Simplifies Deployment:**

➢ Evaluate models using original code & systems

➢ Models run in separate processes

  ➢ Resource isolation

**Common Interface → Simplifies Deployment:**

➢ Evaluate models using original code & systems

➢ Models run in separate processes

    ➢ Resource isolation

    ➢ Scale-out

**Problem:** frameworks optimized for **batch processing** not **latency**

# *Adaptive Batching* to Improve Throughput

➤ Why batching helps:

A single page load may generate many queries

Hardware Acceleration

Helps amortize system overhead

➤ Optimal batch depends on:
  - ➤ hardware configuration
  - ➤ model and framework
  - ➤ system load

**Clipper Solution:**

*be as **slow** as **allowed**…*

➤ Application specifies latency objective

➤ Clipper uses TCP-like tuning algorithm to **increase latency** up to the objective

# Tensor Flow Conv. Net (GPU)

Optimal Batch Size

Latency (ms)

Throughput

Latency Deadline

P99 Latency

Avg. Latency

Throughput (Queries Per Second)

Batch Sizes (Queries)

# Comparison to TensorFlow Serving



**Takeaway**: *Clipper is able to **match the average latency** of TensorFlow Serving while reducing **tail latency (2x)** and **improving throughput (2x)***

# *Approximate Caching to* Reduce Latency

➢ Opportunity for caching

Popular items may be evaluated frequently

➢ Need for **approximation**

Bag-of-Words Model **Images**

High Dimensional and continuous valued queries have low cache hit rate.

**Clipper Solution: *Approximate Caching***

apply *locality sensitive hash functions*



Hash 1

Hash 2

Cache Hit

Cache Miss

Cache Hit

Error

# Clipper Architecture

Applications

**RPC/REST Interface**

Predict ↕ | ↕ Observe

## Clipper

### Correction Layer

Correction Policy

### Model Abstraction Layer

Approximate Caching

Adaptive Batching

RPC ↕ | RPC ↕ | RPC ↕ | RPC ↕

Model Wrapper (MW) | MW | MW | MW

KeystoneML | Caffe | | scikit learn

•••

| Correction Policy | Correction Layer |
|---|---|

**Goal:**

*Maximize **accuracy** through **ensembles, online learning,** and **personalization***

Generalize the **split-model** insight from Velox to achieve:

➢ **robust predictions** by combining multiple models & frameworks

➢ **online learning** and **personalization** by correcting and personalizing **predictions** in response to feedback

# Learning

# Inference

## Velox



Big Data

**Slow** Changing Model

**Fast** Changing User Model

Application

Feedback

Fast Feedback

Slow

# Learning

# Inference

**Slow** Changing Model

## Clipper

**Fast** Changing User Model

Big Data

Feedback

Fast Feedback

Slow

Application

# Correction Policy

Improves prediction **accuaray** by:

➤ Incorporating real-time **feedback**

➤ Managing **personalization**

➤ **Combine** models & **frameworks**

    ➤ enables frameworks to **compete**



**Slow** Changing Model

**Clipper**

**Fast** Changing User Model

# Improved Prediction **Accuracy** (ImageNet)

| System | Model | Error Rate | #Errors |
|--------|-------|-----------:|--------:|
| Caffe | VGG | 13.05% | 6525 |
| Caffe | LeNet | 11.52% | 5760 |
| Caffe | ResNet | 9.02% | 4512 |
| TensorFlow | Inception v3 | 6.18% | 3088 |

sequence of pre-trained state-of-the-art models

# Improved Prediction Accuracy

| System | Model | Top-5 Percent | Errors |
|--------|-------|---------------|--------|
| Caffe | | | 6525 |
| Caffe | | | 5760 |
| Caffe | ResNet | 9.02% | 4512 |
| TensorFlow | Inception v3 | 6.18% | 3088 |
| **Clipper** | **Ensemble** | **5.86%** | **2930** |

**5.2%** relative improvement in prediction accuracy!

# Cost of Ensembles

## Increased Load

➤ *Solutions:*
  ➤ **Caching** and **Batching**
  ➤ **Load-shedding** correction policy can prioritize frameworks

## Stragglers

➤ e.g., framework fails to meet SLO

➤ *Solution:* **Anytime** predictions
  ➤ Correction policy must render predictions with missing inputs
  ➤ e.g., built-in correction policies **substitute expected value**

**Slow** Changing Model

# Clipper

**Fast** Changing User Model

scikit learn

Caffe

?

# Anytime Predictions

# Anytime Predictions

**Fast** Changing User Model

**Slow** Changing Model

$$w_1^{\text{Gogh}} f_{\text{scikit}}(x) + w_2^{\text{Gogh}} \mathbb{E}_X \left[ f_{\text{TF}}(X) \right] + w_3^{\text{Gogh}} f_{\text{Caffe}}(x)$$

Caffe

# Evaluation of Throughput Under Heavy Load



**Takeaway**: *Clipper is able to **gracefully degrade accuracy** to maintain availability under heavy load.*

# Coarsening + Anytime Predictions

$$f_i(x; \theta) \approx$$
$$f_i(z; \theta)$$

$$f_i(x; \theta) \approx$$
$$\mathbb{E}\left[f_i(x; \theta)\right]$$

# Conclusion

Clipper sits between applications and ML frameworks to



**Clipper**

➢ to **simplifying deployment**

➢ **bound latency** and **increase throughput**

➢ and enable **real-time learning** and **personalization**

across **machine learning frameworks**

Big
Data

Training

Model

Query

Decision

Application

Feedback

VELOX

Clipper

NETFLIX

theano Dato Caffe scikit learn
KeystoneML Create TensorFlow mxnet KALDI
VW

# Ongoing & Future Research Directions

➤ Serving and updating RL models

➤ Bandit techniques in correction policies
  ➤ **Collaboration with MSR**

➤ Splitting inference across the cloud and the client to reduce latency and bandwidth requirements

➤ Secure model evaluation on the client (model DRM)